

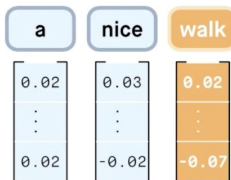


Roadmap



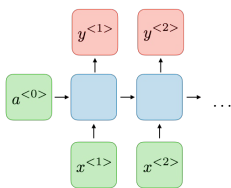
01

Contextualized Embeddings



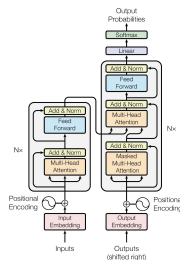
02

From RNNs to Transformers



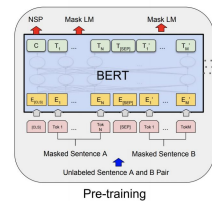
03

Transformers Step-by-Step



04

BERT



05

Distillation and Practical Example





Acknowledgements



Material based on:

- Christoffer Manning's [NLP Lectures at Stanford](#)
- [The Illustrated Transformer](#) by Jay Alammar
- [Slides](#) from Jacob Devlin
- [Self-attention Video](#) from Peltarion

01 / Contextualized Embeddings



Background to Natural Language Processing (NLP)



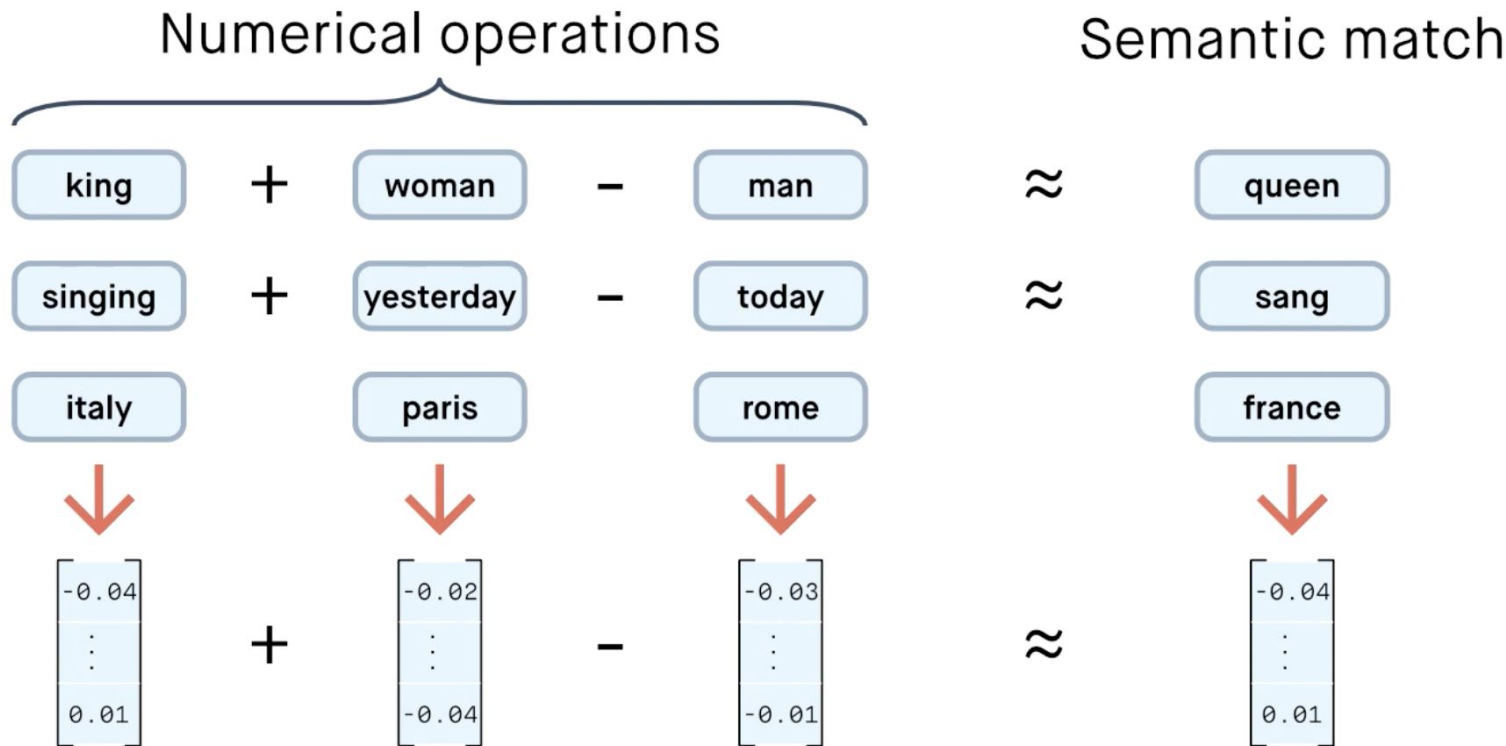
- **Word embeddings** are the basis of NLP
- Popular embeddings like **GloVe** and **Word2Vec** are pre-trained on large text corpuses based on **co-occurrence statistics**
- “A word is characterized by the company it keeps” [Firth, 1957]

best	-	selling	music	artists
↓	↓	↓	↓	↓
-0.11	0.01	-0.01	0.06	-0.02
0.01	0.07	-0.03	0.11	0.00
-0.17	-0.04	0.15	0.05	-0.05
⋮	⋮	⋮	⋮	⋮
0.13	-0.05	0.00	0.14	0.05
-0.13	-0.11	-0.07	-0.12	-0.12
-0.09	-0.25	0.05	-0.04	0.02

[Peltarion, 2020]



Word Embeddings

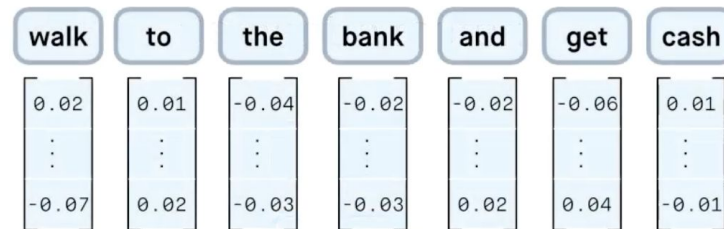
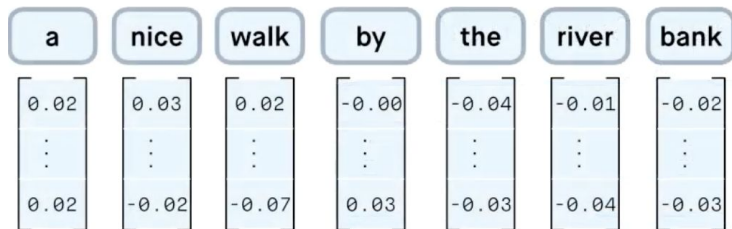




Word Embeddings



Problem: Word embeddings are **context-free**

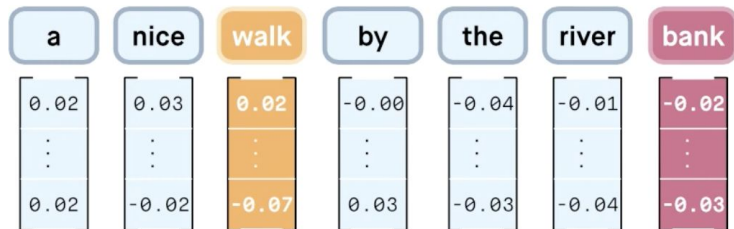




Word Embeddings



Problem: Word embeddings are **context-free**



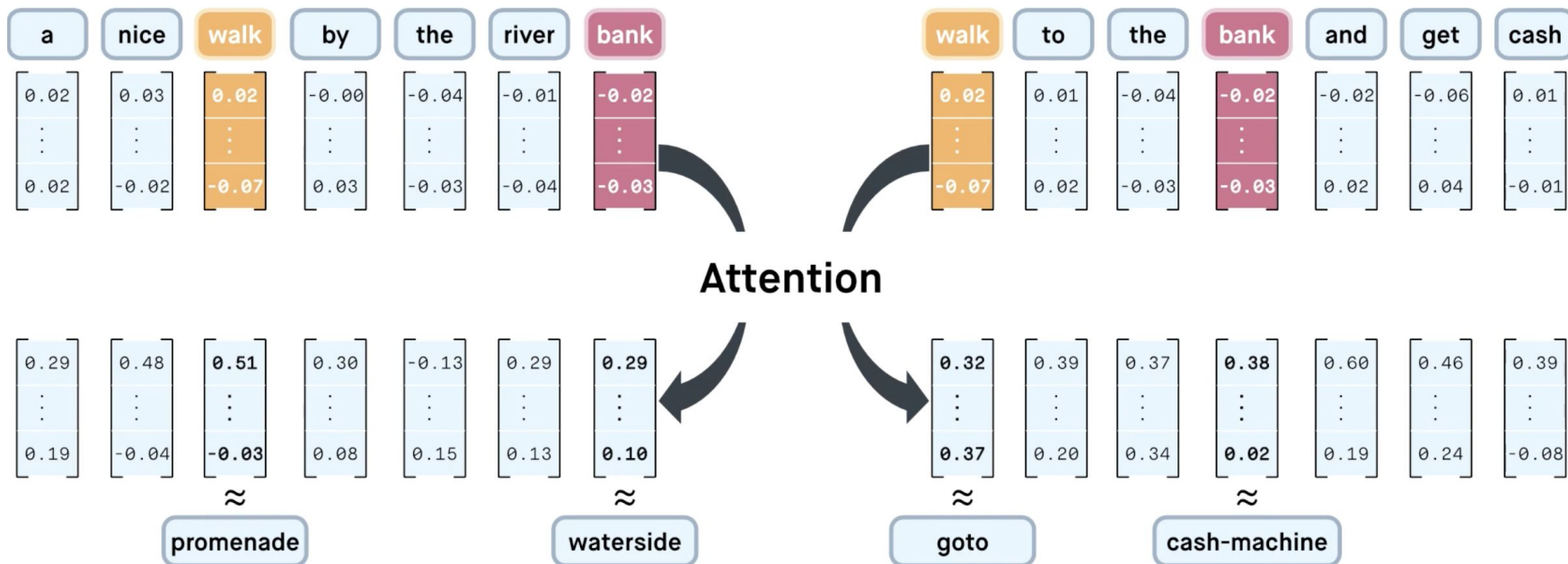


Word Embeddings



Problem: Word embeddings are **context-free**

Solution: Create **contextualized** representation



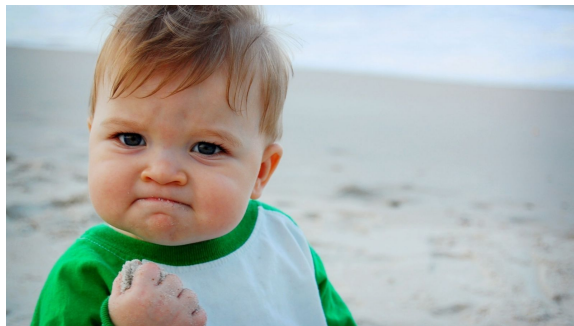
02 / From RNNs to Transformers



Problems with RNNs - Motivation for Transformers



- Sequential computations **prevents parallelization**
- Despite GRUs and LSTMs, RNNs still need attention mechanisms to deal with **long range dependencies**
- Attention gives us access to any state... Maybe we don't need the costly recursion? 🤔
- Then NLP can have deep models, solves our computer vision envy!

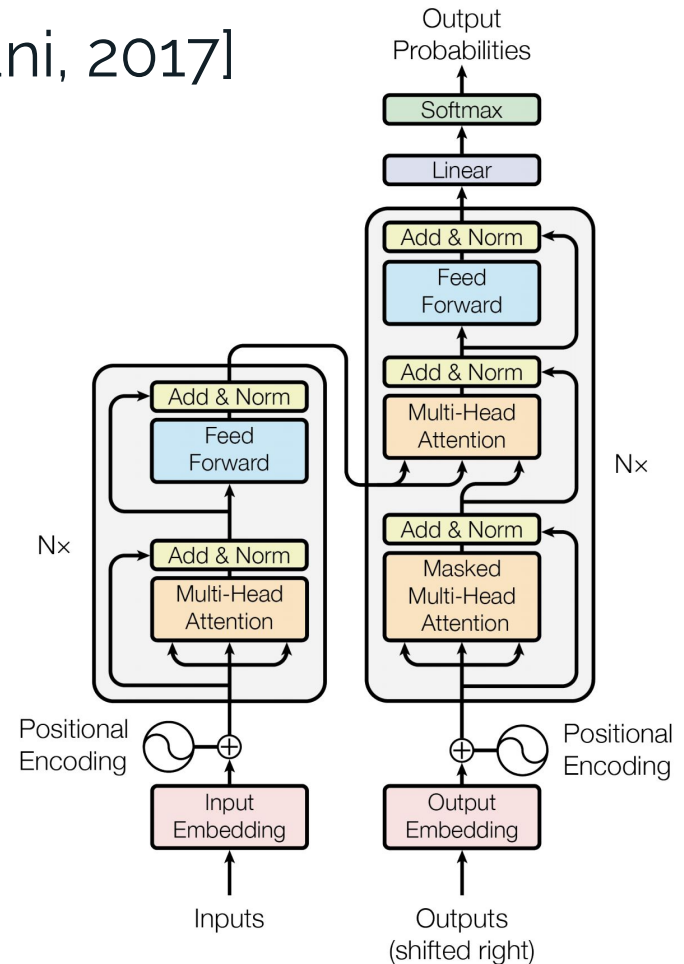




Attention is all you need! [Vaswani, 2017]



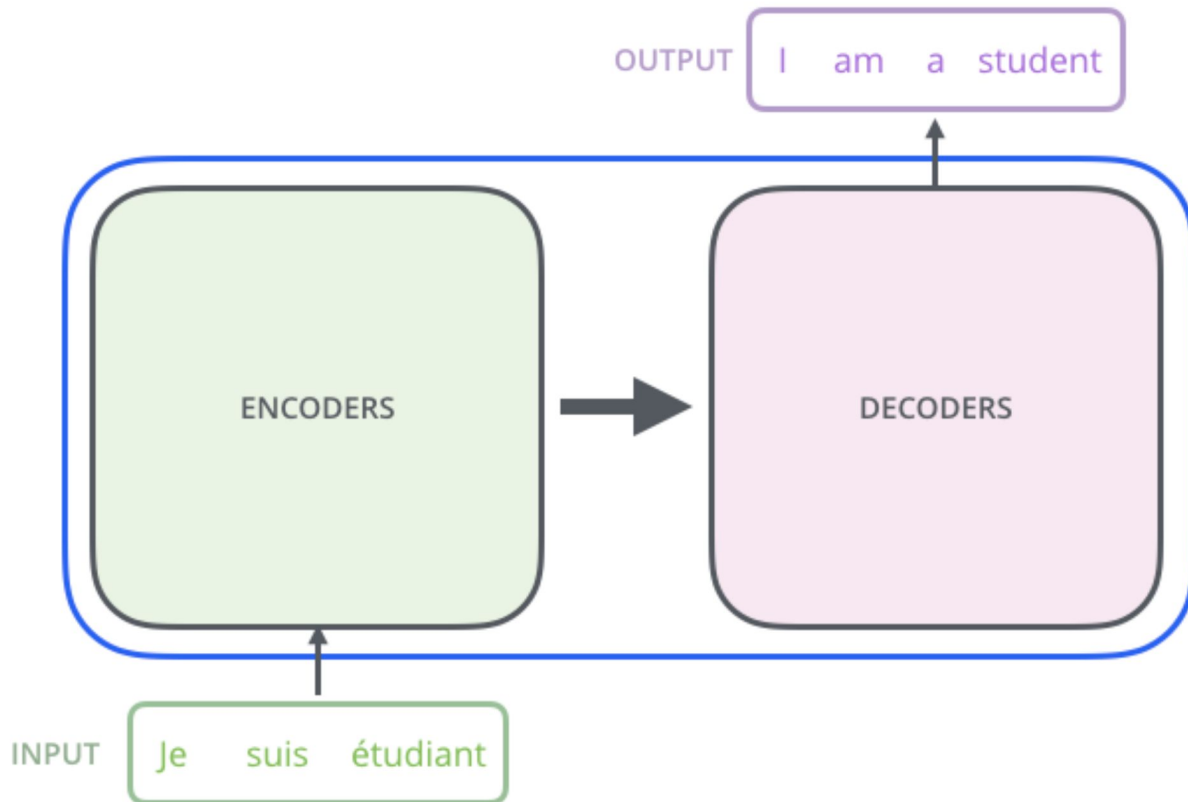
- Sequence-to-sequence model for Machine Translation
- **Encoder-decoder** architecture
- Multi-headed **self-attention**
 - Models context and no locality bias



03 / Transformers Step-by-Step



Understanding the Transformer: Step-by-Step



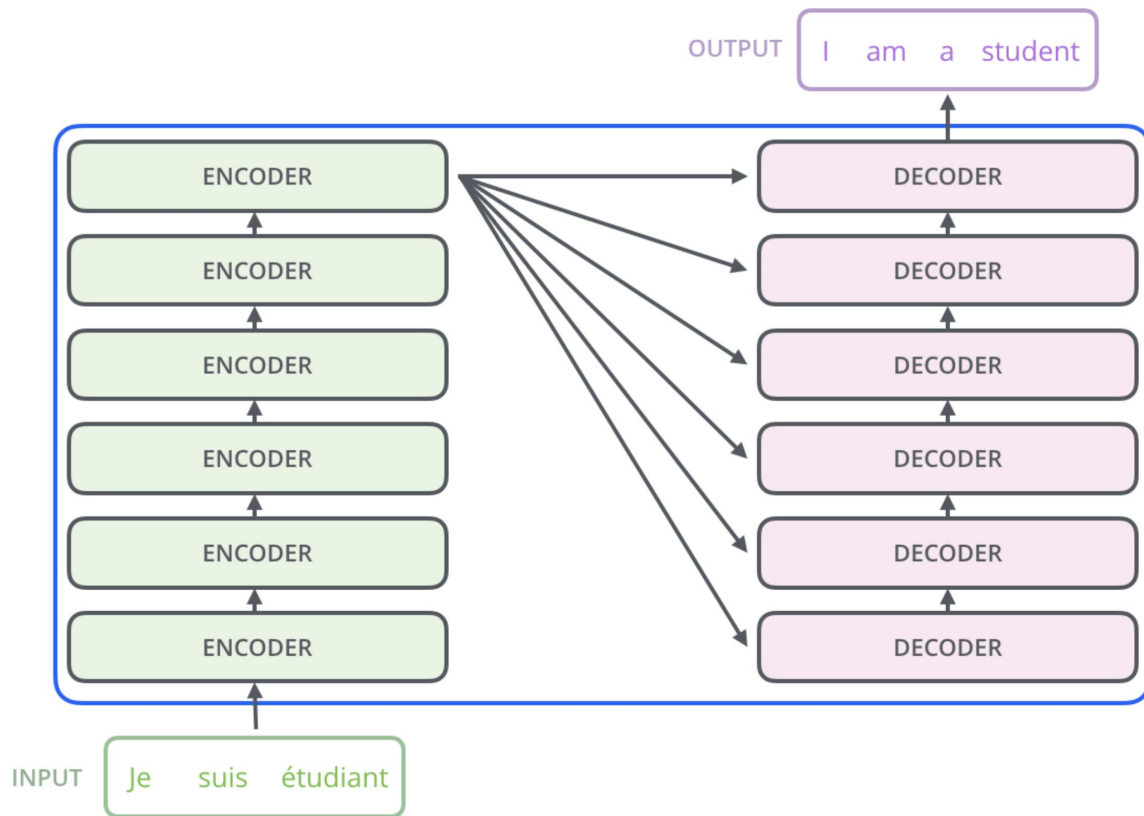


Understanding the Transformer: Step-by-Step



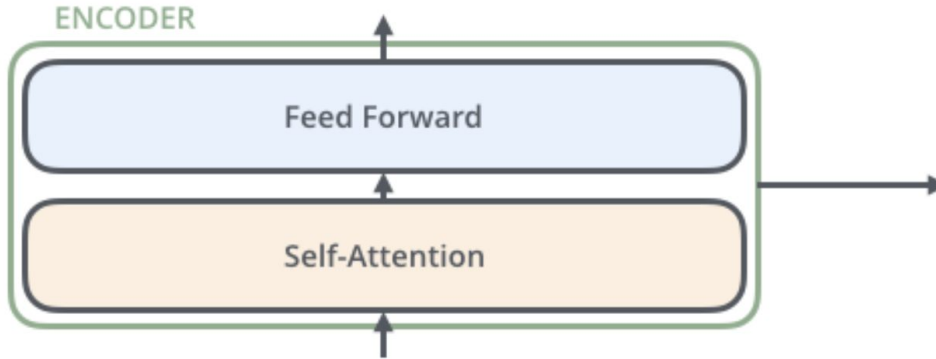
No recursion, instead **stacking encoder** and **decoder** blocks

- Originally: 6 layers
- BERT base: 12 layers
- BERT large: 24 layers
- GPT2-XL: 48 layers
- GPT3: 96 layers

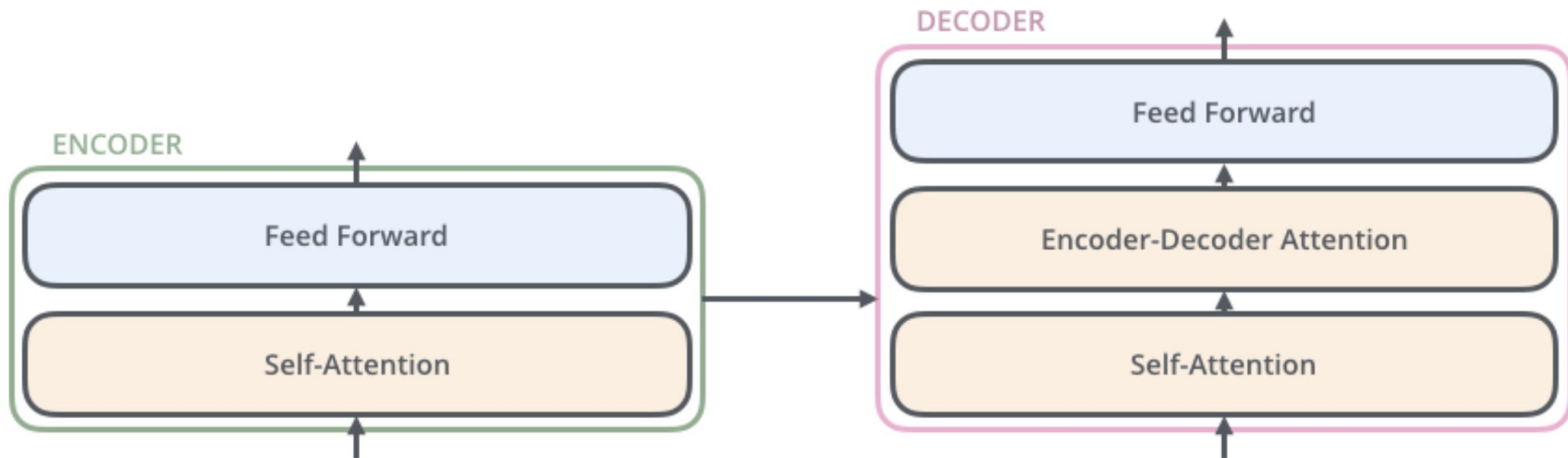




The Encoder Block



The Encoder and Decoder Blocks

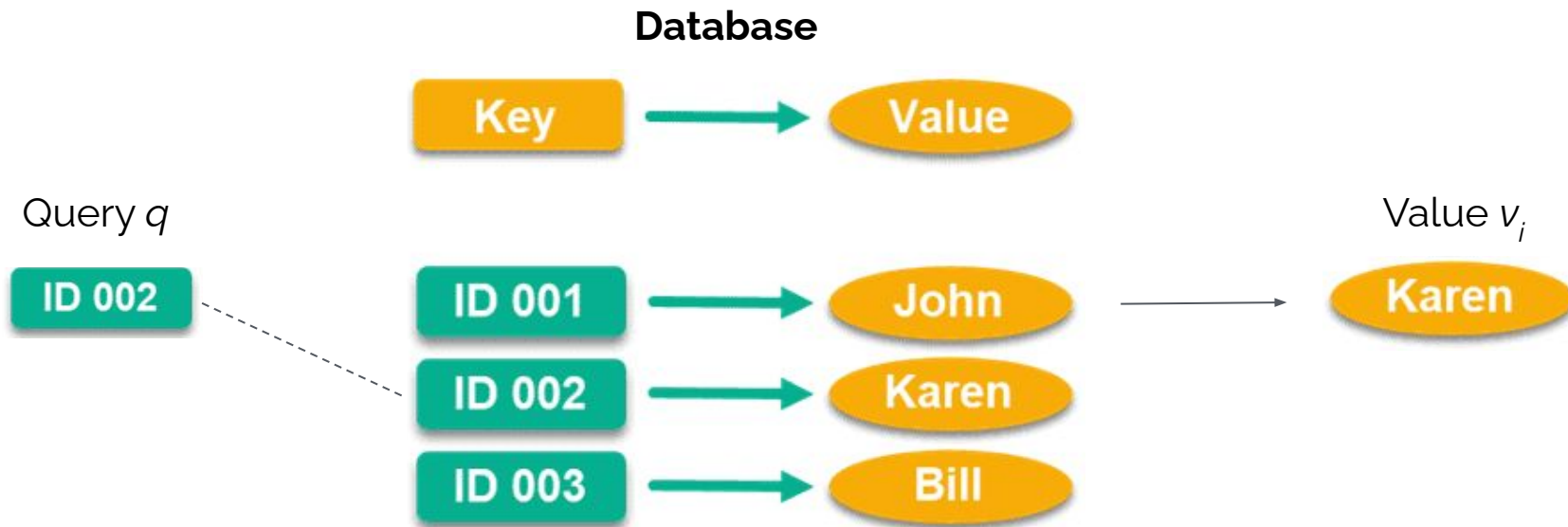




Attention Preliminaries



Mimics the retrieval of a value v_i for a query q based on a key k_i in a database, but in a probabilistic fashion





Dot-Product Attention



- Queries, keys and values are vectors
- Output is a **weighted sum** of the values
- Weights are computed as the **scaled dot-product** (similarity) between the query and the keys

$$\text{Attention}(q, K, V) = \sum_i \text{Similarity}(q, k_i) \cdot v_i = \sum_i \frac{e^{q \cdot k_i / \sqrt{d_k}}}{\sum_j e^{q \cdot k_j / \sqrt{d_k}}} v_i$$

Output is a row-vector

- Can stack multiple queries into a matrix Q

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Output is again a matrix

- Self-attention: Let the word embeddings be the queries, keys and values, i.e. **let the words select each other**



Self-Attention Mechanism

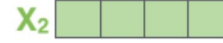
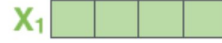


Input

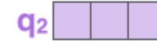
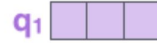
Thinking

Machines

Embedding

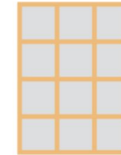
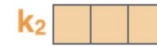
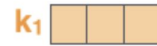


Queries



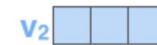
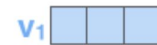
W^Q

Keys



W^K

Values



W^V



Self-Attention Mechanism



Input

Embedding

Queries

Keys

Values

Score

Divide by 8 ($\sqrt{d_k}$)

Softmax

Thinking

x_1

q_1

k_1

v_1

$q_1 \cdot k_1 = 112$

14

0.88

Machines

x_2

q_2

k_2

v_2

$q_1 \cdot k_2 = 96$

12

0.12



Self-Attention Mechanism in Matrix Notation



$$X \times W^Q = Q$$

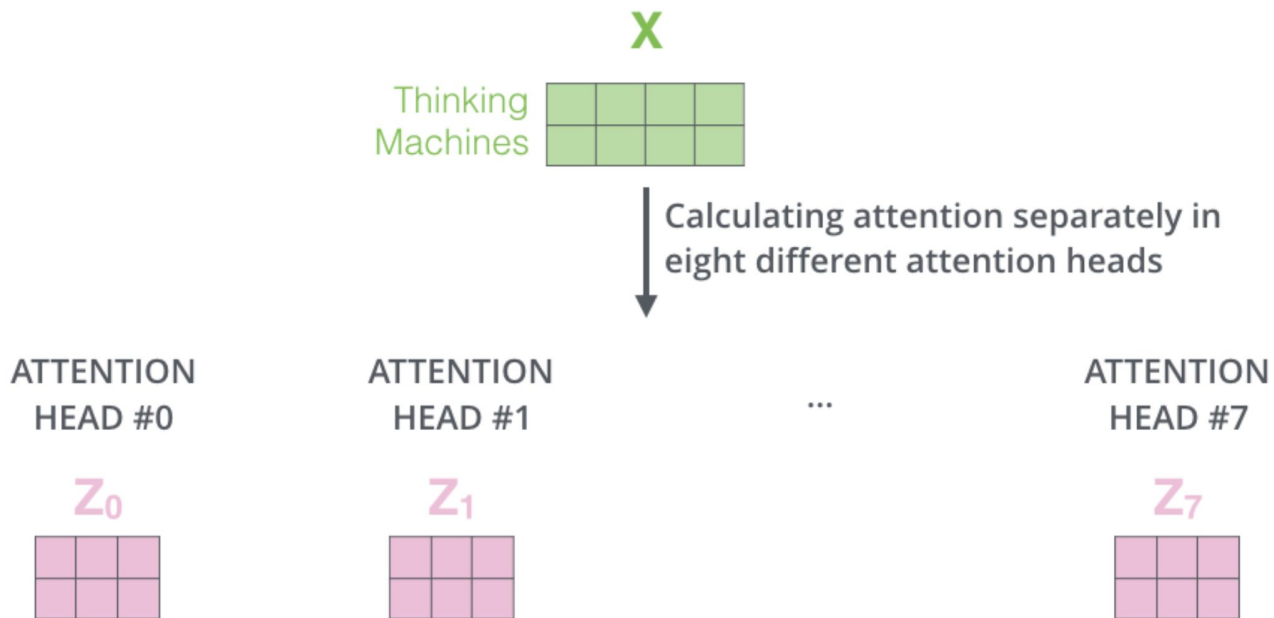
$$X \times W^K = K$$

$$X \times W^V = V$$

$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) \times V = Z$$

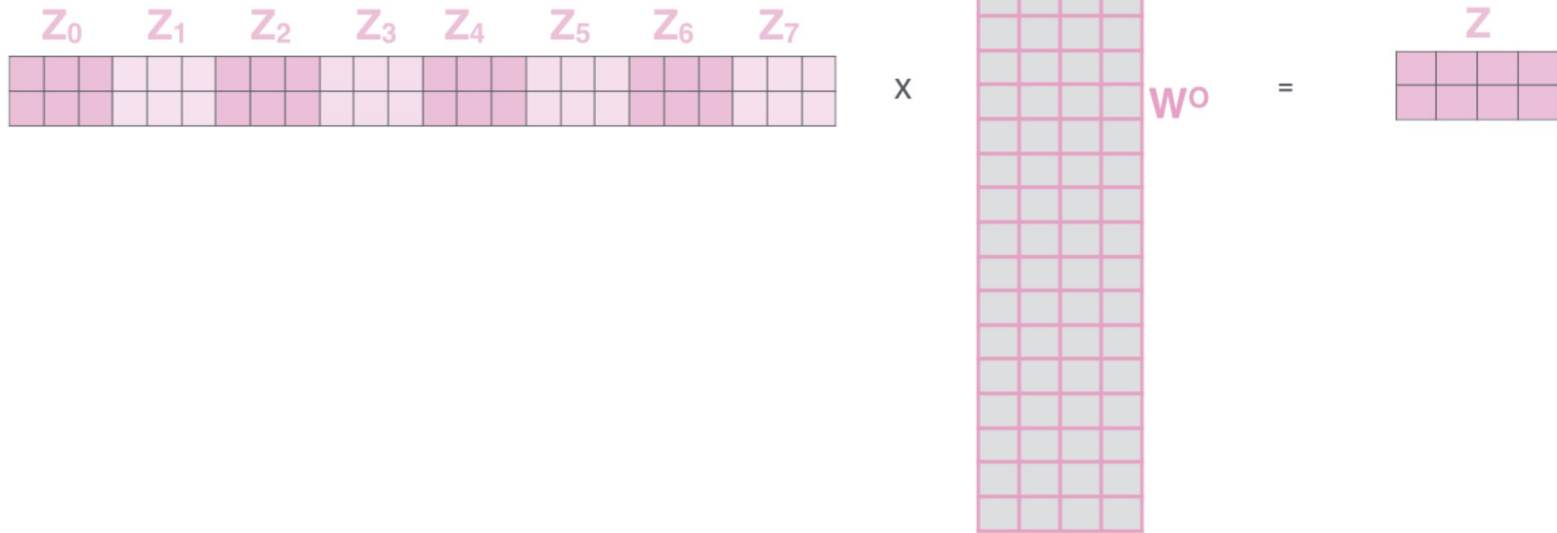


Multi-Headed Self-Attention





Multi-Headed Self-Attention

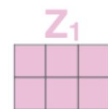
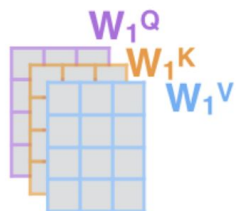
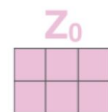
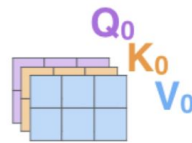
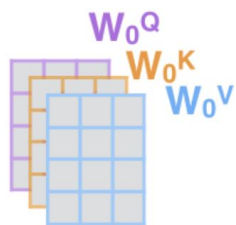




Self-Attention: Putting It All Together



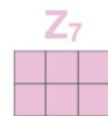
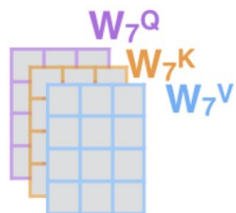
Thinking
Machines



...

...

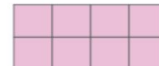
...



W^O

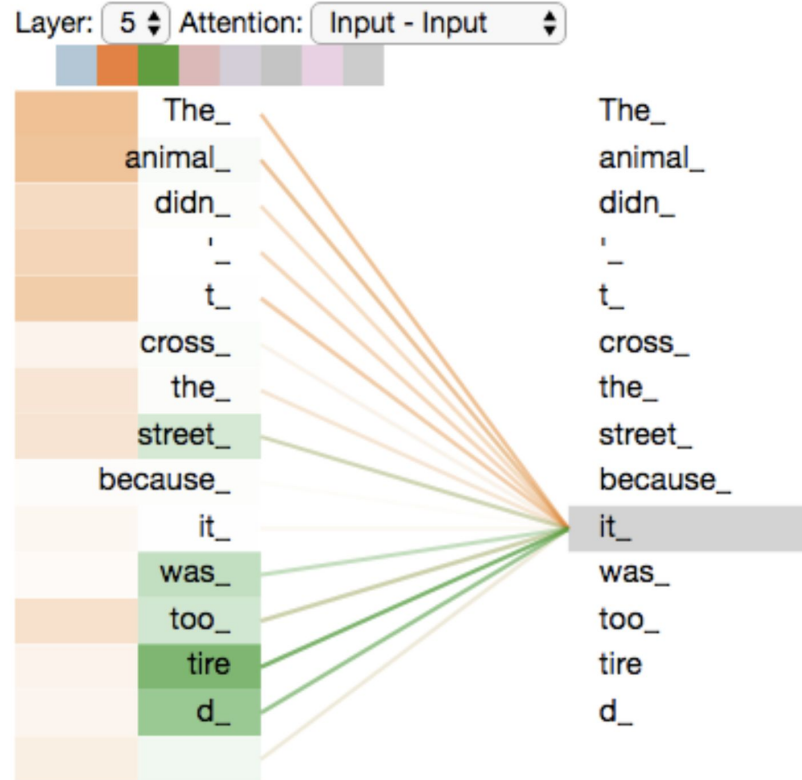


Z





Attention visualized



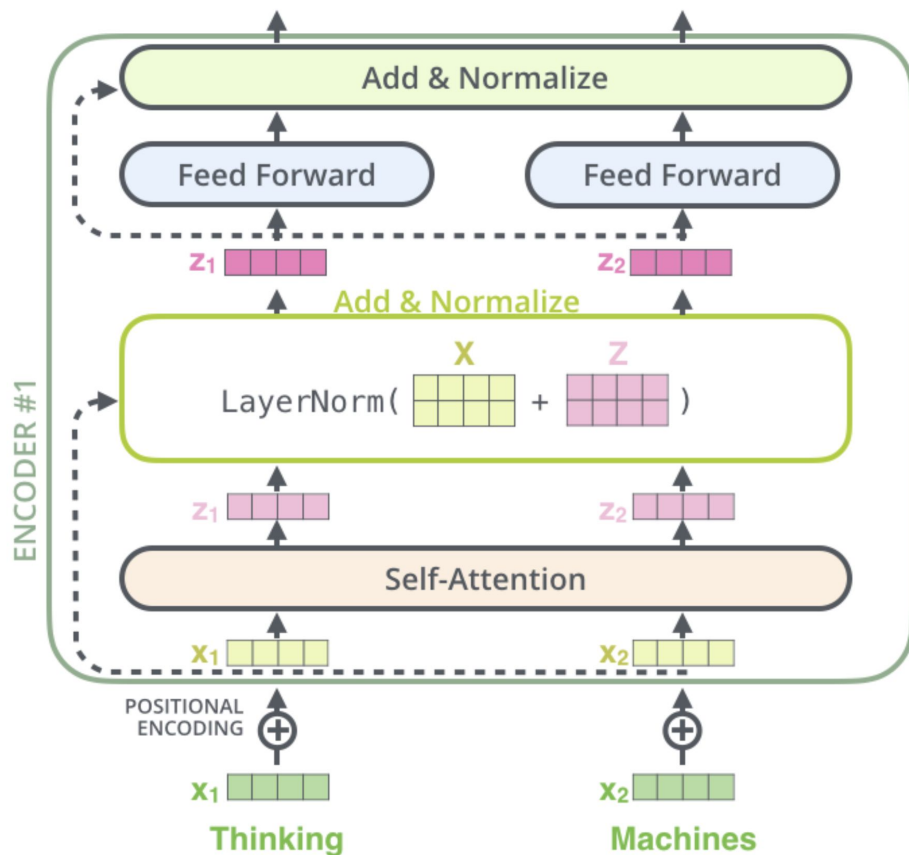


The Full Encoder Block



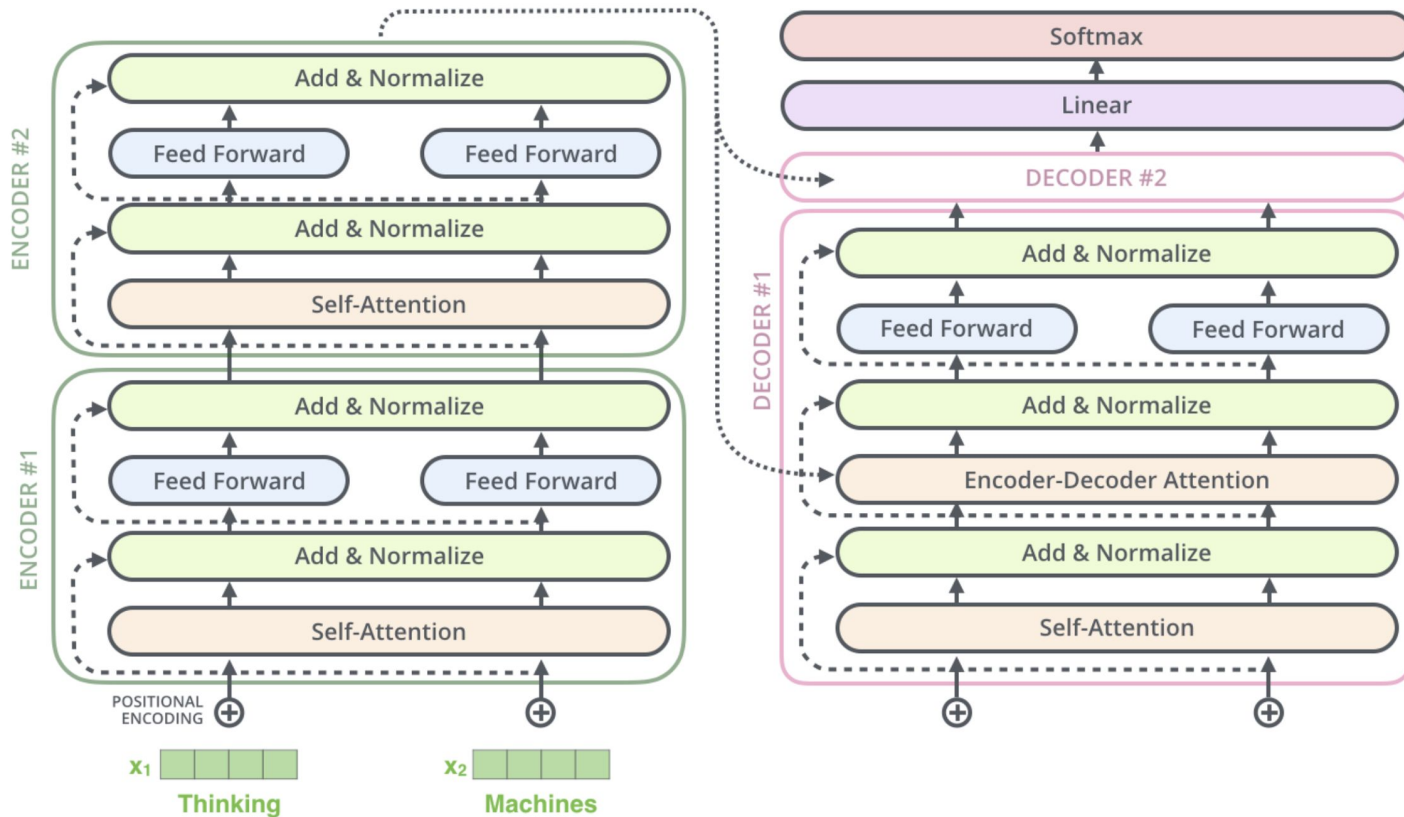
Encoder block consisting of:

- Multi-headed self-attention
- Feedforward NN (FC 2 layers)
- Skip connections
- Layer normalization - Similar to batch normalization but computed over features (words/tokens) for a single sample





Encoder-Decoder Architecture - Small Example





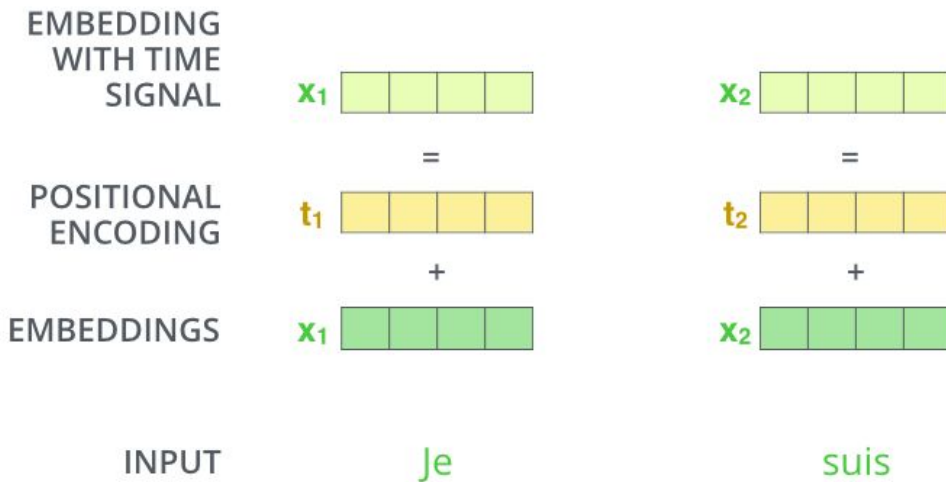
Positional Encodings



- Attention mechanism has no locality bias - **no notion of word order**
- **Add positional encodings** to input embeddings to let model learn relative positioning

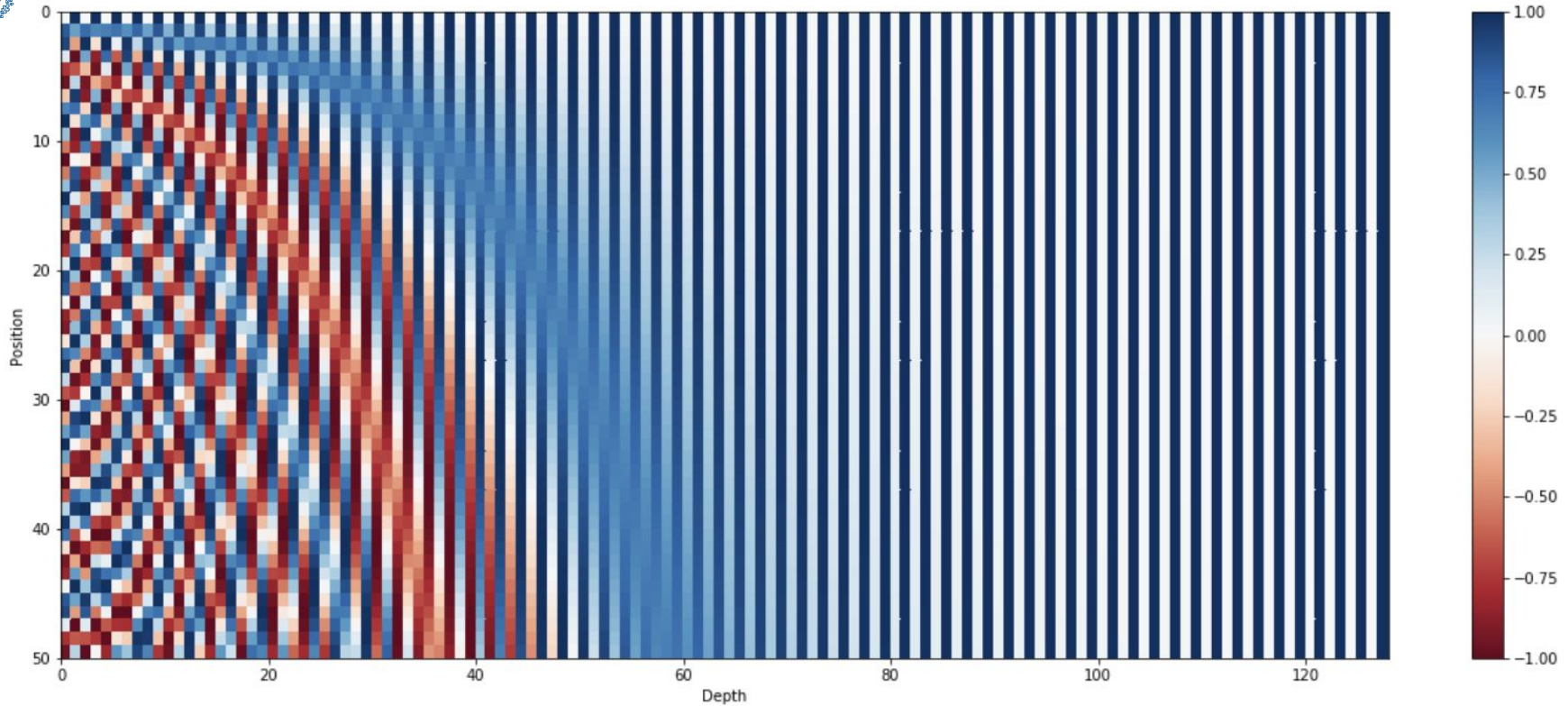
$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

$$PE(pos, 2i + 1) = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$





Positional Encodings



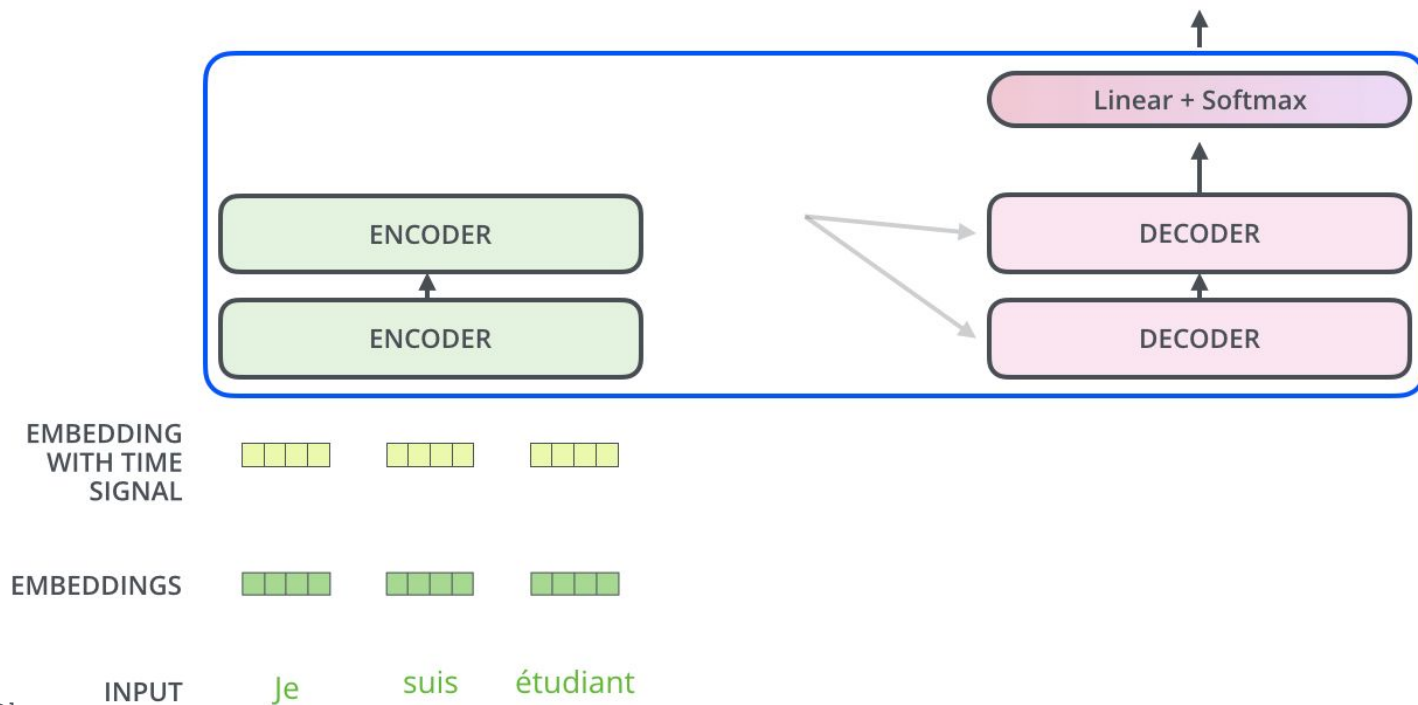


Let's start the encoding!



Decoding time step: 1 2 3 4 5 6

OUTPUT



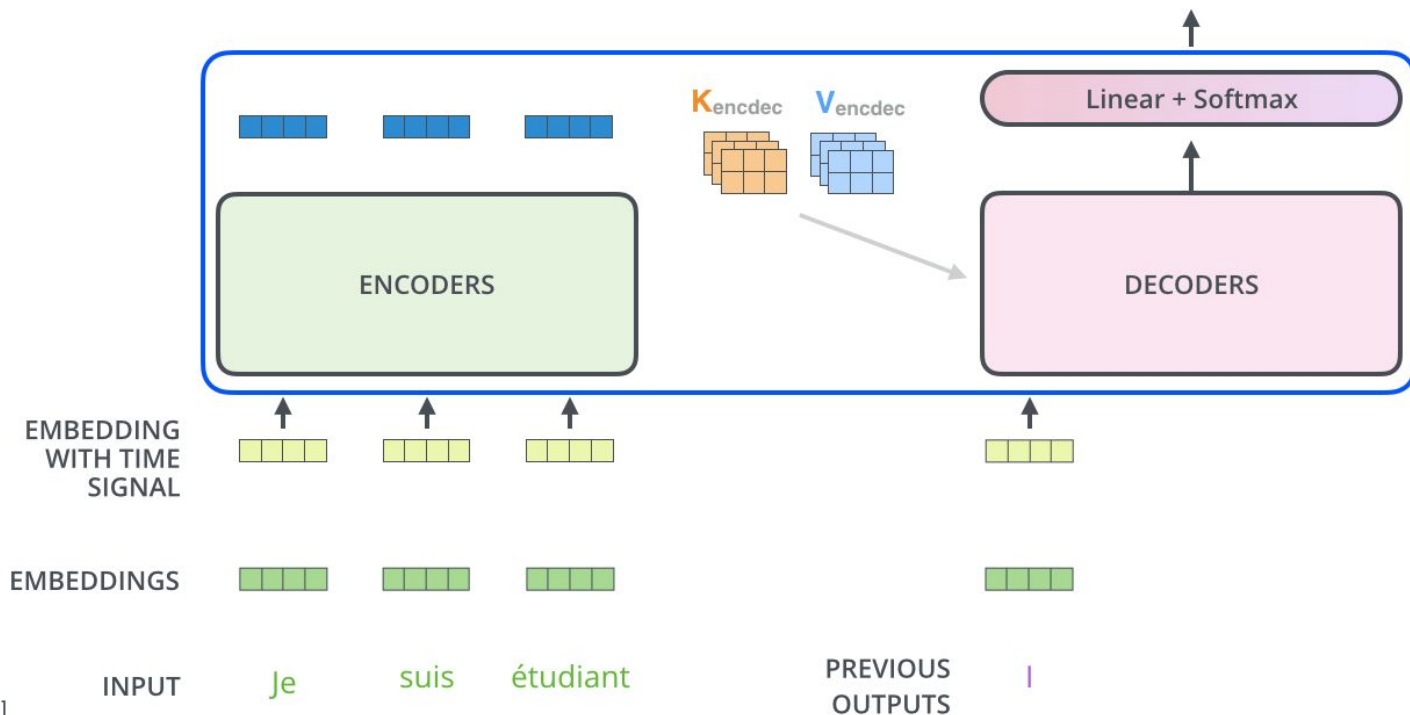


Decoding procedure



Decoding time step: 1 2 3 4 5 6

OUTPUT |

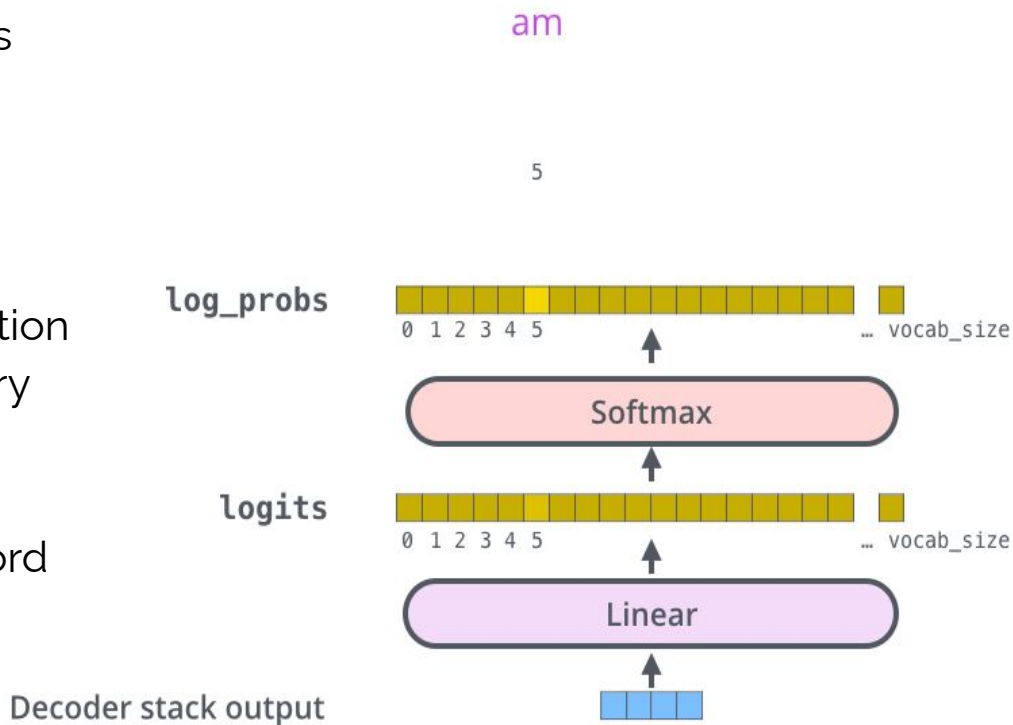




Producing the output text



- The output from the decoder is passed through a final fully connected **linear layer** with a **softmax** activation function
- Produces a probability distribution over the pre-defined vocabulary of output words (tokens)
- **Greedy decoding** picks the word with the highest probability at each time step



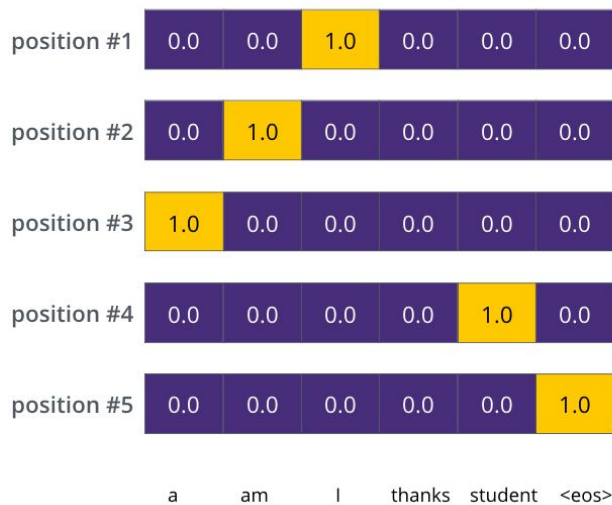


Training Objective



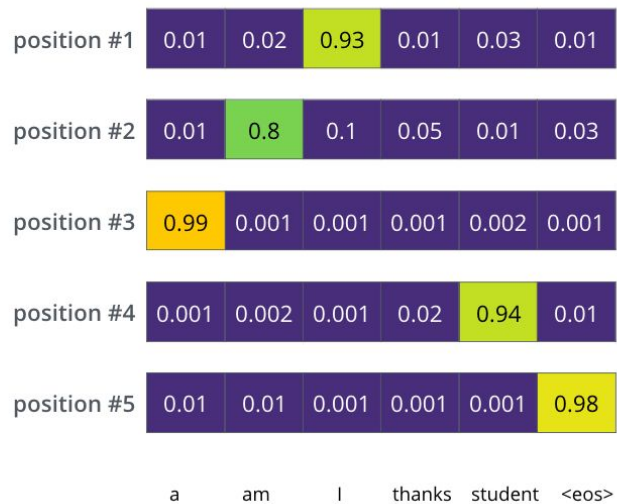
Target Model Outputs

Output Vocabulary: a am I thanks student <eos>



Trained Model Outputs

Output Vocabulary: a am I thanks student <eos>





Complexity Comparison



Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$

04 / BERT



BERT



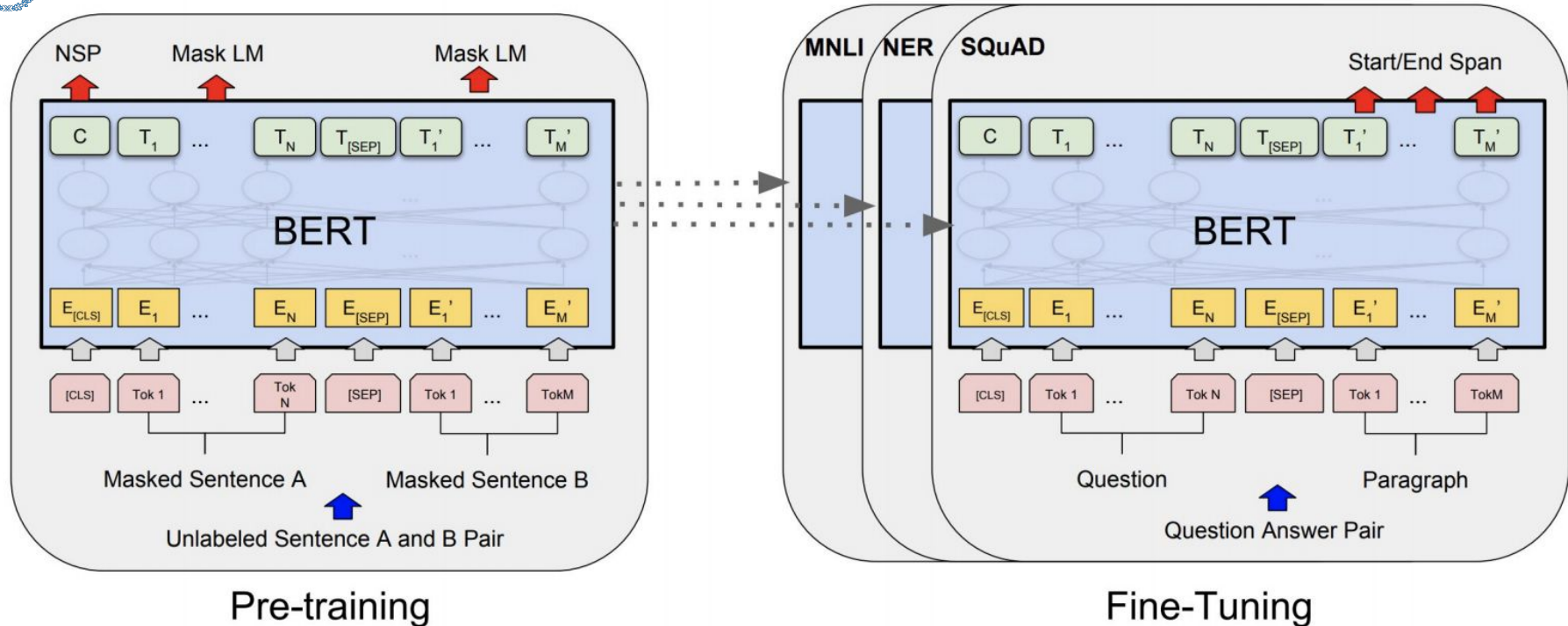
Bidirectional **E**ncoder **R**epresentations from **T**ransformers

- Self-supervised **pre-training** of Transformers encoder for **language understanding**
- **Fine-tuning** for specific downstream task





BERT Training Procedure





BERT Training Objectives



Masked Language Modelling

the man went to the [MASK] to buy a [MASK] of milk

store gallon

↑ ↑

Next Sentence Prediction

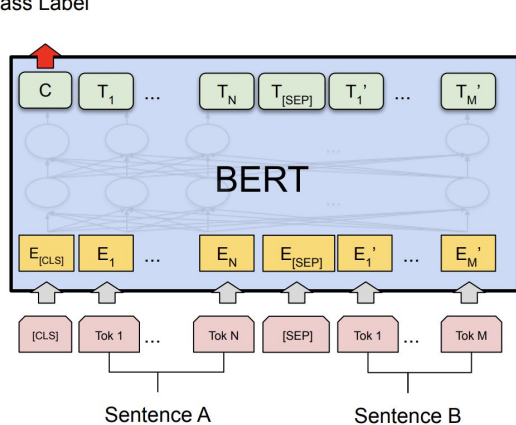
Sentence A = The man went to the store.
Sentence B = He bought a gallon of milk.
Label = IsNextSentence

Sentence A = The man went to the store.
Sentence B = Penguins are flightless.
Label = NotNextSentence

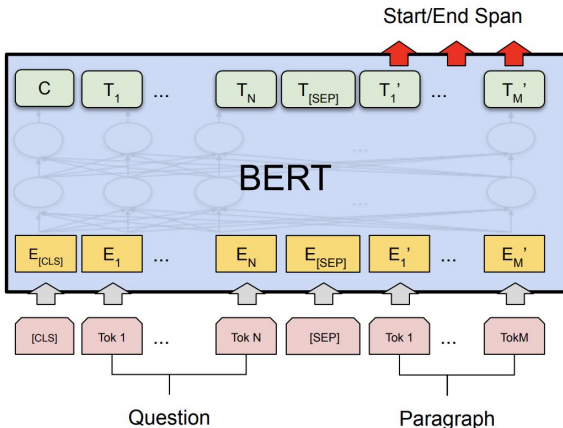
BERT Fine-Tuning Examples



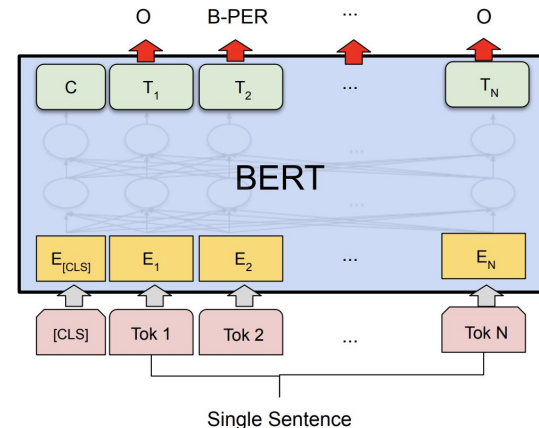
Class Label



**Sentence
Classification**



**Question
Answering**



**Named Entity
Recognition**



Exploring the Limits of Transfer Learning (T5)



- Scaling up **models size** and amount of **training data** helps a lot
- Best model is 11B (!!) parameters
- Exact **pre-training objective** (MLM, NSP, corruption) doesn't matter too much
- SuperGLUE benchmark:

Rank	Name	Model	URL	Score	BoolQ	CB	COPA	MultiRC	ReCoRD	RTE	WiC	WSC	AX-b	AX-g
1	SuperGLUE Human Baselines	SuperGLUE Human Baselines		89.8	89.0	95.8/98.9	100.0	81.8/51.9	91.7/91.3	93.6	80.0	100.0	76.6	99.3/99.7
+ 2	T5 Team - Google	T5		89.3	91.2	93.9/96.8	94.8	88.1/63.3	94.1/93.4	92.5	76.9	93.8	65.6	92.7/91.9
+ 3	Huawei Noah's Ark Lab	NEZHA-Plus		86.7	87.8	94.4/96.0	93.6	84.6/55.1	90.1/89.6	89.1	74.6	93.2	58.0	87.1/74.4
+ 4	Alibaba PAI&ICBU	PAI Albert		86.1	88.1	92.4/96.4	91.8	84.6/54.7	89.0/88.3	88.8	74.1	93.2	75.6	98.3/99.2
+ 5	Tencent Jarvis Lab	RoBERTa (ensemble)		85.9	88.2	92.5/95.6	90.8	84.4/53.4	91.5/91.0	87.9	74.1	91.8	57.6	89.3/75.6
6	Zhuiyi Technology	RoBERTa-mtl-adv		85.7	87.1	92.4/95.6	91.2	85.1/54.3	91.7/91.3	88.1	72.1	91.8	58.5	91.0/78.1
7	Facebook AI	RoBERTa		84.6	87.1	90.5/95.2	90.6	84.4/52.5	90.6/90.0	88.2	69.9	89.0	57.9	91.0/78.1

05 / Practical Examples



BERT in low-latency production settings



GOOGLE \ TECH \ ARTIFICIAL INTELLIGENCE \

Google is improving 10 percent of searches by understanding language context

Say hello to BERT

By [Dieter Bohn](#) | [@backlon](#) | Oct 25, 2019, 3:01am EDT

Bing says it has been applying BERT since April

The natural language processing capabilities are now applied to all Bing queries globally.

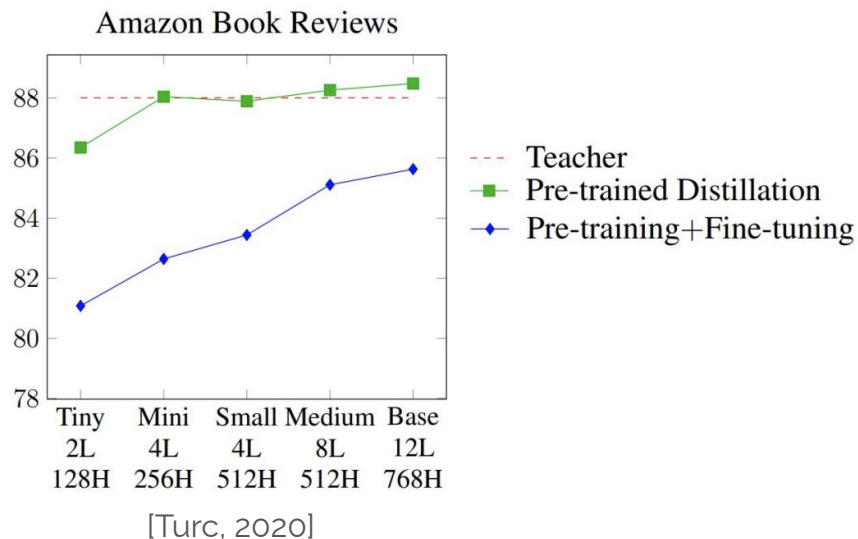
[George Nguyen](#) on November 19, 2019 at 1:38 pm



Distillation



- Modern pre-trained language models are **huge** and very **computationally expensive**
- How are these companies applying them to low-latency applications?
- Distillation!
 - Train SOTA **teacher model** (pre-training + fine-tuning)
 - Train smaller **student model** that **mimics** the teacher's output on a large dataset on unlabeled data
- Why does it work so well?





Transformers in TensorFlow using HuggingFace 🤗



- The **HuggingFace Library** contains a majority of the recent pre-trained State-of-the-art NLP models, as well as over 4 000 community uploaded models
- Works with both **TensorFlow** and **PyTorch**

The screenshot shows the HuggingFace website interface. At the top, there is the HuggingFace logo (a yellow emoji) and the text 'HUGGING FACE'. Below this, there is a link 'Back to home' and the heading 'All Models and checkpoints'. A yellow box contains the text: 'Also check out our list of [Community contributors](#) 🏆 and [Organizations](#) 🌐.' Below this is a search bar with the text 'Search models...', a 'Tags: All' dropdown, and a 'Sort: Most downloads' dropdown. A list of model names is displayed, each with a star icon to its right:

- bert-base-uncased
- deepset/bert-large-uncased-whole-word-masking-squad2
- distilbert-base-uncased
- dccuchile/bert-base-spanish-wwm-cased
- microsoft/xprophetnet-large-wiki100-cased-xglue-ntg
- deepset/roberta-base-squad2
- jplu/tf-xlm-roberta-base
- cl-tohoku/bert-base-japanese-whole-word-masking
- distilroberta-base
- bert-base-cased
- xlm-roberta-base



Transformers in TensorFlow using HuggingFace 🤗



```
from transformers import BertTokenizerFast, TFBertForSequenceClassification
from datasets import load_dataset
import tensorflow as tf

dataset = load_dataset("imdb").shuffle()
tokenizer = BertTokenizerFast.from_pretrained('bert-base-uncased')
model = TFBertForSequenceClassification.from_pretrained('bert-base-uncased', num_labels=2)

train_encodings = tokenizer(dataset['train']['text'], truncation=True, padding=True)
train_dataset = tf.data.Dataset.from_tensor_slices((dict(train_encodings), dataset['train']['label']))
val_dataset = ... // Analogously

optimizer = tf.keras.optimizers.Adam(learning_rate=5e-5)
model.compile(optimizer=optimizer, loss=model.compute_loss)
model.fit(train_dataset.batch(16), epochs=3, batch_size=16)

model.evaluate(val_dataset.batch(16), verbose=0)
```



Transformers in TensorFlow using HuggingFace 🤗



```
from transformers import BertTokenizerFast, TFBertForSequenceClassification
from datasets import load_dataset
import tensorflow as tf

dataset = load_dataset("imdb").shuffle()
tokenizer = BertTokenizerFast.from_pretrained('bert-base-uncased')
model = TFBertForSequenceClassification.from_pretrained('bert-base-uncased', num_labels=2)

train_encodings = tokenizer(dataset['train']['text'], truncation=True, padding=True)
train_dataset = tf.data.Dataset.from_tensor_slices((dict(train_encodings), dataset['train']['label']))
val_dataset = ... // Analogously

optimizer = tf.keras.optimizers.Adam(learning_rate=5e-5)
model.compile(optimizer=optimizer, loss=model.compute_loss)
model.fit(train_dataset.batch(16), epochs=3, batch_size=16)

model.evaluate(val_dataset.batch(16), verbose=0)
```



Transformers in TensorFlow using HuggingFace 🤗



```
from transformers import BertTokenizerFast, TFBertForSequenceClassification
from datasets import load_dataset
import tensorflow as tf
```

```
dataset = load_dataset("imdb").shuffle()
tokenizer = BertTokenizerFast.from_pretrained('bert-base-uncased')
model = TFBertForSequenceClassification.from_pretrained('bert-base-uncased', num_labels=2)
```

```
train_encodings = tokenizer(dataset['train']['text'], truncation=True, padding=True)
train_dataset = tf.data.Dataset.from_tensor_slices((dict(train_encodings), dataset['train']['label']))
val_dataset = ... // Analogously
```

```
optimizer = tf.keras.optimizers.Adam(learning_rate=5e-5)
model.compile(optimizer=optimizer, loss=model.compute_loss)
model.fit(train_dataset.batch(16), epochs=3, batch_size=16)
```

```
model.evaluate(val_dataset.batch(16), verbose=0)
```




Transformers in TensorFlow using HuggingFace 🤗



```
from transformers import BertTokenizerFast, TFBertForSequenceClassification
from datasets import load_dataset
import tensorflow as tf
```

```
dataset = load_dataset("imdb").shuffle()
tokenizer = BertTokenizerFast.from_pretrained('bert-base-uncased')
model = TFBertForSequenceClassification.from_pretrained('bert-base-uncased', num_labels=2)
```

```
train_encodings = tokenizer(dataset['train']['text'], truncation=True, padding=True)
train_dataset = tf.data.Dataset.from_tensor_slices((dict(train_encodings), dataset['train']['label']))
val_dataset = ... // Analogously
```

```
optimizer = tf.keras.optimizers.Adam(learning_rate=5e-5)
model.compile(optimizer=optimizer, loss=model.compute_loss)
model.fit(train_dataset.batch(16), epochs=3, batch_size=16)
```

```
model.evaluate(val_dataset.batch(16), verbose=0)
```

06 / Wrap Up



Summary



- Transformers have blown other architectures out of the water for NLP
- Get rid of recurrence and rely on **self-attention**
- NLP pre-training using **Masked Language Modelling**
- Most recent improvements using **larger models** and **more data**
- **Distillation** can make model serving and inference more tractable



Questions?

/November 25, 2020

— **Thanks**

Karl Fredrik Erliksson

PhD Candidate, KTH and Peltarion
kferl@kth.se